

1 A framework for predicting tissue-specific 2 effects of rare genetic variants

3 Farhan N. Damani¹, Yungil Kim¹, Xin Li², Emily K. Tsang³, Joe R. Davis⁴,
4 Colby Chiang⁵, Zachary Zappala⁴, Benjamin J. Strober⁶, Alexandra J.
5 Scott⁵, Ira M. Hall⁵, GTEx Consortium, Stephen B. Montgomery^{2,4}, and
6 Alexis Battle¹

7 ¹Department of Computer Science, Johns Hopkins University, Baltimore, MD.

8 ²Department of Pathology, Stanford University, Stanford, CA.

9 ³Biomedical Informatics Program, Stanford University, Stanford, CA.

10 ⁴Department of Genetics, Stanford University, Stanford, CA.

11 ⁵McDonnell Genome Institute, Washington University School of Medicine, St. Louis,
12 MO.

13 ⁶Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD.

14 Corresponding author:

15 Alexis Battle¹

16 Email address: ajbattle@cs.jhu.edu

17 ABSTRACT

18 Despite the abundance of rare genetic variants—variants carried by less than one
19 percent of the population—in human genomes, the impact of these variants on specific
20 tissues has been largely uncharacterized. Population-level test statistics, while effective
21 in understanding the impact of common variants—variants carried by at least five
22 percent of the population, have had limited success in characterizing the effect of
23 rare variants mainly due to limited statistical power. In addition, the effect of each
24 rare variant can vary greatly between specific tissues. This heterogeneity coupled
25 with limited sample sizes and a lack of known disease-causing rare variants makes
26 predicting tissue-specific cellular consequences of rare variants a difficult task. To
27 make these predictions, we propose a new method called SPEER (SPecific tissuE
28 variant Effect predictoR): a hierarchical Bayesian model that uses transfer learning,
29 allowing separate predictions in each tissue while flexibly sharing signal across tissues
30 to improve power. Our probabilistic model capitalizes on a growing body of rich
31 epigenetic annotations to inform the consequences of a variant in specific tissues.
32 These annotations are integrated with tissue-specific RNA expression levels and
33 common variants. We show our method improves prediction accuracy in simulations
34 and in genomic data from the Genotype-Tissue Expression (GTEx) project.

INTRODUCTION

Recent advances in genomic technologies provide us with a unique opportunity to study the contribution of genetic variation to disease risk. Genome-wide association studies (GWAS) have been largely successful over the past decade in identifying statistical associations between common genetic variants—those carried by at least five percent of the population, and complex traits and diseases including height, diabetes and heart disease. However, these statistical techniques do not generalize well to analyzing rare variants—variants carried by less than one percent of the population—due to low sample size (Uricchio et al., 2016). Because rare variants have been shown to be implicated in disease risk and shown to be potentially more deleterious than common variants (Tennessen et al., 2012; Nelson et al., 2012), developing methods that can effectively characterize these variants remains essential.

Several tools have been developed to understand the functional consequences of rare variants. Kircher et al. (2014) developed CADD, a supervised learning approach that used functional annotations of the genome to predict deleteriousness. Quang et al. (2015) built on the success of CADD with a deep learning approach also for predicting deleteriousness. Li et al. (2016) introduced RIVER, an unsupervised learning method that integrates genomic annotations with gene expression data from the same individual to prioritize deleterious variants. They showed that genomic annotations are enriched for variants nearby genes with extreme expression levels. Building on this knowledge, RIVER used gene expression outliers—samples with extreme over or under expression—across diverse tissues to prioritize deleterious variants. They were better able to identify deleterious variants with global effects compared to models that exclusively used genomic annotations.

While these methods have made significant strides in understanding the global impact of genetic variants, their usefulness in understanding the tissue-specific consequences of genetic variants is somewhat limited. Recent work by Backenroth et al. (2016) integrated tissue-specific regulatory elements with GWAS summary statistics in order to understand these effects. Despite providing unique insights about the sharing of genetic variants within known physiological tissue groups Aguet et al. (2016), these methods do not apply to rare variant analysis due to a lack of known pathogenic tissue-specific rare variants and a scarcity of samples.

Transfer learning, a framework that allows sharing of knowledge across learning tasks, has been shown to be effective in low-resource settings with complex structure (Thrun, 1996; McCallum et al., 1998). In the hierarchical Bayes framework, parameters for each task are dependent on each other through a Bayesian prior (Raina et al., 2006). Here we propose SPEER (SPecific tissuE variant Effect predictor), a hierarchical Bayes model that uses transfer learning to predict the tissue-specific functional consequences of rare variants. Each task here translates to understanding the effects of rare variants in a specific tissue and the sharing across tissues captures global effects. By using transfer learning to share information across tissues, SPEER learns reliable parameters and prioritizes rare variants in a tissue-specific manner. SPEER has three parts. First, a per-

80 sample component models the effect of both genomic annotations and gene expression
81 on the presence of rare regulatory variation. Second, a tissue-specific component models
82 the influence of genomic annotations on individual tissues. Third, a global component
83 models the shared impact of genomic annotations across tissues.

84

85 We apply our method to simulated data and data from the Genotype-Tissue Expression
86 (GTEx) project and show that SPEER performs better than state-of-the-art baselines.
87 The methods developed in this paper are available at <https://github.com/farhand7/speer>.

88

89

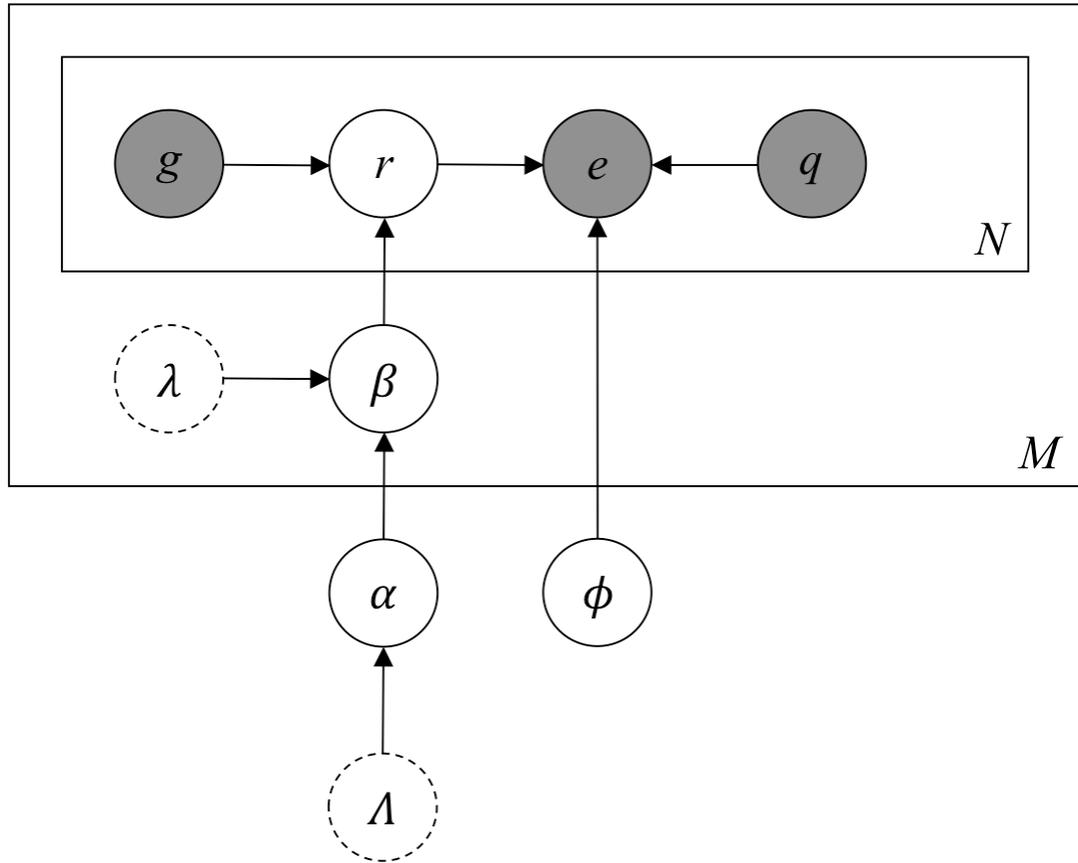


Figure 1. Graphical representation of our model. The outer plate represents tissues, while the inner plate represents individuals and genes within a tissue. Shaded circles represent observed variables; white circles represent hidden variables; dotted edged circles represent hyperparameters.

90 METHODOLOGY

91 SPEER is a probabilistic model for inferring the functional consequences of rare variants
 92 in M individual tissues. For each tissue c , we have N_c samples, each representing a
 93 single individual for a single gene. For each sample i within tissue c , \mathbf{X} posits that
 94 the presence of a rare regulatory variant r_{ci} can be inferred by integrating measured
 95 tissue-specific gene expression e_{ci} , significant common variants q_{ci} nearby sample i , and
 96 genomic annotations g_{ci} describing the rare variants nearby sample i , which is a function
 97 of both tissue-specific $\{\beta_c, \lambda_c\}$ and shared tissue parameters $\{\alpha, \Lambda\}$. The graphical
 98 model is shown in Figure 1.

99
 100 SPEER infers the presence of a rare regulatory variant nearby a sample by optimizing a
 101 joint objective function. The objective has three components: a global component, a
 102 tissue-specific component, and a sample-level component.

103
 104

$$\begin{aligned}
\log p(e, g, r, q, \beta, \lambda, \alpha, \Lambda, \phi) = & \underbrace{\log p(\alpha|\Lambda)}_{\text{(A) global component}} + \underbrace{\sum_{c=1}^M \left(\sum_{j=1}^L \log p(\beta_{cj}|\alpha_j, \lambda_c) \right)}_{\text{(B) tissue-specific component}} \\
& + \underbrace{\sum_{i=1}^{N_c} \log \sum_{r_{ci}}^S p(e_{ci}|r_{ci}, q_{ci}, \phi) p(r_{ci}|g_{ci}, \beta_c)}_{\text{(C) per-sample component}}
\end{aligned} \tag{1}$$

105

106 **Per-sample component.** Each individual by gene sample is assumed to belong to one
107 of S latent groups (functional variant classes). The random variable $r_{ci} \in \{1, \dots, S\}$
108 encodes functional variant class membership. We infer the membership of each sample
109 by integrating genomic annotations, tissue-specific gene expression, and significant
110 common variants. $g_{ci} \in \mathbb{R}^L$ is a vector of L genomic annotations describing the set of
111 rare variants nearby sample i , and $\beta_c \in \mathbb{R}^L$ is a vector of L weights. Formally, we model
112 the effects of g_{ci} on r_{ci} as:

$$r_{ci}|g_{ci}, \beta_c \sim \text{Bern}(\psi)$$

$$\psi = \frac{1}{1 + e^{-\beta_c^T g_{ci}}}$$

113 We expect functional variants to cause disruption at a cellular level potentially evident by
114 individual molecular phenotypes. Similar to Li et al., we hypothesize that extreme gene
115 expression levels can inform effects of rare variants even at low frequencies. Therefore,
116 we use tissue-specific gene expression outliers denoted by $e_{ci} \in \{0, 1\}$, which identifies
117 the outlier status of sample i within tissue c . We compute outliers by evaluating whether
118 the absolute z-score of a sample's gene expression is greater than a predefined threshold.
119 $q_{ci} \in \{0, 1\}$ denotes the presence of a significant common variant nearby the gene in
120 sample i . Together we model the effects of r_{ci}, q_{ci} , on e_{ci} as:

$$e_{ci}|r_{ci}, q_{ci}, \phi \sim \text{NoisyOr}(\phi)$$

121 ϕ controls the rate of functional rare variants to expression outliers and is the same
122 across tissues.

123

124 **Tissue-specific component.** Genomic annotations g_{ci} are assumed to inform both
125 global and tissue-specific effects of genetic variants. For each tissue c , $\beta_c \in \mathbb{R}^L$ is a
126 random variable that deviates from the global effects parameter $\alpha \in \mathbb{R}^L$ with a tissue-
127 specific transfer factor $\lambda_c \in \mathbb{R}$. λ_c is shared across features. For the j th feature, we
128 have:

$$\beta_{cj}|\alpha_j, \lambda_c \sim \mathcal{N}(\alpha_j, \lambda_c^{-1})$$

129 We exclusively model transferable effects between tissues, not between tissue-specific
130 features. This allows our model to scale well with a large number of annotations.

131

132 **Global component.** The shared tissue level captures global effects across tissues.
 133 For the j th feature, the global genomic annotations coefficients $\alpha_j \in \mathbb{R}^L$ is distributed as
 134 $\alpha_j | \Lambda \sim \mathcal{N}(\vec{0}, \Lambda^{-1})$.

135

136 Learning

137 We want to learn the parameters of our model $\Theta = \{\beta_{1:M_G}, \phi, \alpha\}$ and our hyperparameters
 138 $\{\lambda_{1:M}, \Lambda\}$.

139

140 We use the empirical Bayes bootstrap estimation procedure described in Efron and
 141 Tibshirani (1994) to estimate the transfer factors $\{\lambda_{1:M}, \Lambda\}$. Let $\delta_{j,c} = \beta_j^c - \alpha_j$. For i
 142 $= \{1, \dots, K\}$ randomly sampled with replacement datasets, we compute the maximum
 143 likelihood estimation (with regularization) for β_c and α . With these estimates, we
 144 compute the empirical variance of $\delta_{j,c}$ across K datasets:

$$\lambda_c^{-1} = \frac{\sum_{i=1}^K \sum_{j=1}^L (\beta_{c,j}^{(i)} - \alpha_j^{(i)})^2}{(K-1)L}$$

145 After estimating our hyperparameters, we compute MAP estimates of Θ by optimizing
 146 the log of the joint distribution in Eq. (1) with respect to Θ . Because latent variables
 147 make optimization non-convex, we use expectation maximization (EM) to maximize the
 148 observed data log likelihood.

149

150 **Expectation step.** We compute the posterior distribution over the set of latent variables
 151 r by conditioning on the observed data and our model parameters. Assuming each
 152 sample is i.i.d, compute:

153

$$q_{ci}(r_{ci}) = p(r_{ci} = 1 | e_{ci}, g_{ci}, q_{ci}, \beta_c, \lambda_c, \alpha, \Lambda, \phi) = \frac{p(r_{ci} = 1 | g_{ci}, \beta_c) p(e_{ci} | r_{ci} = 1, q_{ci}, \phi)}{\sum_{r_{ci}} p(r_{ci} | g_{ci}, \beta_c) p(e_{ci} | r_{ci}, q_{ci}, \phi)} \quad (2)$$

154 **Maximization step.** The expectation of the complete data log likelihood with respect
 155 to $p(\mathbf{r} | \dots)$ is:

156

157

$$\arg \max_{\beta_{1:M_L}, \alpha, \phi} \log p(\alpha | \Lambda^{-1}) + \sum_{c=1}^M \left(\sum_{j=1}^L \log p(\beta_{c,j} | \alpha_j, \lambda_c^{-1}) \right) + \sum_{i=1}^{N_c} \sum_{z_{ci}} q(r_{ci}) \log [p(r_{ci} | g_{ci}, \beta_c) p(e_{ci} | r_{ci}, q_{ci}, \phi)] \quad (3)$$

158 We use blocked coordinate gradient descent to estimate β_c and α , iterating between
 159 updating $\alpha_j = \frac{\sum_{c=1}^M \lambda_c \beta_{c,j}}{\Lambda + \sum_{c=1}^M \lambda_c}$ and $\beta_{c,j}^{t+1} = \beta_{c,j}^t - \nabla f(\beta_{c,j}^t, \alpha_j^t, q_{ci}, g_{ci})$, where $\nabla f = \frac{\partial f}{\partial \beta_{c,j}^t} =$

160 $-\lambda_c(\beta_{cj}^t - \alpha_j^t) + \sum_{i=1}^{N_c} -g_{cij}(q_{ci}(r_{ci}) - h(\beta_c, g_{ci}))$ where h is the inverse logit function.
161 ϕ is updated using a NoisyOR MAP estimation procedure with soft assignments to \mathbf{r} as
162 weights.

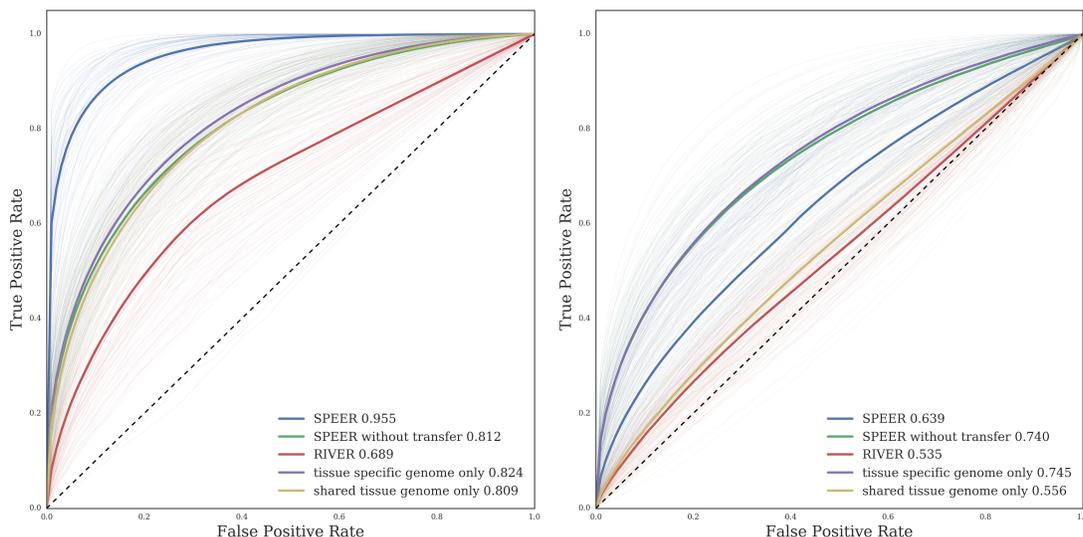


Figure 2. Tissue-specific receiver operating characteristic (ROC) curve averaged across five tissue groups using the stronger effects parameter setting evaluated in the tied tissue (A) and independent tissue (B) simulations. The darker lines represent average ROC curves across 75 simulated runs. Area under the curve (AUC) scores are reported in the legend. There are four benchmarks described here: SPEER w/o transfer was trained on the same data as SPEER but assumes parameter independence; RIVER integrates genomic annotations with shared tissue expression outlier status in an unsupervised setting; shared tissue genome only is a supervised model trained on exclusively genomic annotations using shared tissue expression outlier status as labels; tissue-specific genome only is also a supervised model trained on exclusively genomic annotations using tissue-specific expression outlier status as labels.

164 **Simulation Results.**

165 To highlight the intuition behind SPEER, we performed two simulations: one involving
 166 tied tissues and the other involving independent tissues. The tied tissue simulation used
 167 transfer learning to generate data. The independent tissue simulation generated data
 168 for each tissue independently. Because none of the other approaches considered here
 169 include common variants, we excluded q in order to evaluate the usefulness of tissue
 170 sharing in simulation. Therefore, tissue-specific gene expression is only conditioned
 171 on r , so we used a categorical distribution with parameter ϕ to model this dependency
 172 and used a Beta prior on ϕ with hyperparameters $\mu_{r_{ci}}$ and $\sigma_{r_{ci}}^2$ to generate the rate of
 173 functional variants to expression outliers. Formally, for sample i within tissue c we have:

$$e_{ci}|r_{ci}, \phi \sim \text{Cat}(\phi)$$

$$\phi_{r_{ci}} \sim \text{Beta}(\mu_{r_{ci}}, \sigma_{r_{ci}}^2)$$

174 We re-parameterized the Beta distribution using a mean and variance (Ferrari and Cribari-
 175 Neto, 2004) to allow for better interpretability of the parameter settings described in
 176 our simulation. The simulations were crafted to mimic scenarios with strong effects
 177 from genomic annotations coupled with noisy gene expression data. Besides the caveat
 178 described above, the simulated data for the tied tissue setting was generated by sampling
 179 from the joint distribution assumed by SPEER, as described in Eq. (1). The independent
 180 tissue setting followed a similar procedure except each β_{c_j} was sampled independently
 from $\mathcal{N}(0, \lambda_c)$. Tables 1 and 2 describe three scenarios that were tested.

Table 1. Tied tissue simulation.

Parameter	stronger effects	equal effects	weaker effects
Λ	0.01	0.01	0.1
λ_c	{2, ..., 6}	{2, ..., 6}	{2, ..., 6}
$\phi_{e z=0} \sim \text{Beta}(\mu, \sigma^2)$	(0.4, 1e-4)	(0.3, 1e-4)	(0.4, 1e-4)
$\phi_{e z=1} \sim \text{Beta}(\mu, \sigma^2)$	(0.6, 1e-4)	(0.7, 1e-4)	(0.6, 1e-4)

Table 2. Independent tissue simulation.

Parameter	stronger effects	equal effects	weaker effects
λ_c	0.01	0.01	0.1
$\phi_{e z=0} \sim \text{Beta}(\mu, \sigma^2)$	(0.4, 1e-4)	(0.3, 1e-4)	(0.4, 1e-4)
$\phi_{e z=1} \sim \text{Beta}(\mu, \sigma^2)$	(0.6, 1e-4)	(0.7, 1e-4)	(0.6, 1e-4)

181

182

183 For each setting, we measured the simulation uncertainty by performing each experi-
 184 ment 75 times. The stronger effects scenario underlined strong influence of genomic
 185 annotations coupled with noisy expression labels. The tied tissue simulation (Table 1)
 186 highlighted genomic annotations with strong functional effects combined with correlated
 187 influences across tissues. The independent tissue simulation showed similarly strong
 188 functional consequences from genomic annotations but independent influences across
 189 tissues. SPEER performed significantly better than all baselines at predicting held-out
 190 tissue-specific labels in the tied simulation (Fig. 2A). Even with limited training data,
 191 SPEER provided a significant performance boost (Fig. 5). In the independent tissue
 192 simulation, SPEER performed worse than the other two tissue-specific models—SPEER
 193 without transfer and tissue specific genome only (Fig. 2B). In this simulation, the data
 194 generation process was independent for each tissue, so encouraging tissue similarity
 195 would rightly hurt performance.

196

197 The equal effects scenario mimicked strong influence of genomic annotations simi-
 198 lar to the previous simulation but with highly predictive gene expression labels. We
 199 observed a significant boost in AUC scores across all benchmarks when using expression
 200 data that is more predictive of regulatory status (Fig. 6A). SPEER scores remained
 201 highly predictive of the regulatory status of rare variants when switching to expression
 202 data with a stronger signal (AUC of 0.988 vs 0.955). We observed significant perfor-
 203 mance boosts in all other models using this parameter setting, implying that highly

204 predictive expression data might be critical to the performance of these other models.

205

206 The weaker effects scenario highlights weaker influence of genomic annotations. In the
207 tied case, we observed a lower AUC score for SPEER compared to the stronger effects
208 scenario (Fig. 7A). Despite a lower AUC score, the predictive performance of SPEER
209 remains significantly better than all other models. In the independent tissue simulation
210 (Fig. 7B), we predictably observed SPEER without transfer performing better than
211 SPEER. Given the weaker influence of genomic annotations, the general performance
212 across all models is worse.

213 **Results from GTEx data.**

214 We applied our method to data from the Genotype-Tissue Expression (GTEx V6p)
215 project. We included whole genome sequence data from 113 donors with European
216 ancestry and 5574 RNA-sequence samples from 27 tissues. We defined a rare variant
217 using a minor allele frequency (MAF) below 1% within the GTEx cohort and within the
218 European panel of the 1000 Genomes project (Consortium, 2015). We restricted our
219 analysis to rare single nucleotide variants (SNVs), which are polymorphisms occurring
220 at specific positions in the genome. We generated a set of genomic features describing
221 each rare SNV. This included describing the location of the rare variant with respect
222 to regulatory elements, the conservation status, and summary statistics from genome
223 only variant predictor tools including CADD and DANN. We also separately generated
224 a set of binary tissue-specific annotations that described whether each rare SNV was
225 present in any of the cell-type specific promoter or enhancer regions from ROADMAP
226 Epigenomics and ENCODE projects (Consortium, 2012; Kundaje et al., 2015) using
227 summary statistics from ChromImpute developed by Ernst and Kellis (2015). We then
228 mapped these annotations to one of the 27 GTEx tissues considered here. We then
229 aggregated all rare SNVs within 10 kb of the transcription start site (TSS) to generate
230 gene-level summary statistics by computing the maximum of each annotation across all
231 nearby rare SNVs. Next, we removed technical and environmental confounders from
232 each tissue's gene expression using PEER estimates (Stegle et al., 2012). We then com-
233 puted gene expression outliers using the z-score across all subjects and genes for each
234 tissue. We refer interested readers to Li et al. (2016) for a complete description of the
235 genomic annotations used, the processing of RNA-expression data and the subsequent
236 gene expression outlier calls. Finally, we identified the top significant common variant
237 nearby each gene using the methods described in Aguet et al. (2016) and used this data
238 to denote the presence of a significant common variant for each sample.

239

240 We measured the sensitivity of our results to the threshold used to call tissue-specific
241 expression outliers in the supplement (Fig. 8). The remaining results used a 1.5 z-score
242 threshold. Because single tissue gene expression outliers are too noisy, we identified
243 clusters of tissues that shared similar patterns of gene expression. We used five tissue
244 groups—brain, digestive, epithelial, artery and fats together, and muscles—as input
245 to our model. We used prior experiments to choose tissue groups by evaluating the
246 pairwise-similarity between individual tissues. A list of tissues in each tissue group is
247 available in the supplement.

248

249 Allele-specific expression is known to present strong evidence of a causal cis-regulatory
 250 effect, which often arises from a non-coding variant (Zhang et al., 2009; Yan et al., 2002).
 251 Because the majority of the rare variants in GTEx are non-coding and heterozygous,
 252 measuring tissue-specific allele-specific expression allowed us to evaluate SPEER at
 253 prioritizing functional variants. We measured allelic imbalance as a function of reference
 254 and alternate allele expression read counts, which is computed using an allelic ratio =
 255 $\left| \frac{ref}{ref-alt} \right| - 0.5$. Higher values here imply greater allelic imbalance.

256
 257 We computed the statistical association between SPEER’s predictions and measured
 258 tissue-specific allelic imbalance for all genes in each tissue using Fisher’s Exact Test
 259 and observed significantly greater predictive power using SPEER compared to all bench-
 260 marks (Fig. 3A). We also measured the effect for each tissue individually (Fig. 9). In
 261 addition, we investigated the SPEER posteriors for samples with strong allelic imbalance
 262 (defined using 90 percentile cut-off) and limited to at least one model having a
 263 posterior greater than 0.5 (Fig. 3B). Among samples with observed allelic imbalance,
 264 SPEER identified 120 samples with all predictions greater than 0.85. Genome only
 265 tissue-specific model identified 3 samples; and the shared tissue genome only model
 identified 2 samples.

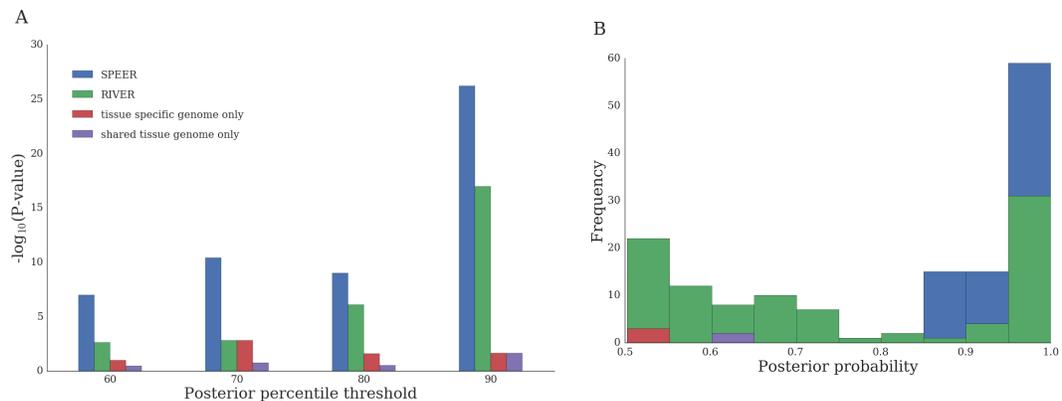


Figure 3. A) Using SPEER scores to predict tissue-specific allelic imbalance. Allelic imbalance was defined by the 90th percentile of allelic ratios. A deleterious SPEER score was defined using four percentile thresholds. We computed p-values for each of the four settings using Fisher’s exact test and compared our results to two benchmarks. B) Histogram of SPEER scores for samples with allelic imbalance limited to samples with at least one of the four models having a posterior greater than 0.5.

266
 267

268 Comparing SPEER to RIVER.

269 SPEER is a probabilistic model that uses transfer learning to infer the *tissue-specific*
 270 regulatory impact of each rare SNV. RIVER is a general method to infer the *global*
 271 regulatory impact of each rare SNV across diverse tissues. SPEER integrates genomic
 272 annotations with tissue-specific expression labels across M tissues. RIVER integrates
 273 genomic annotations with a shared tissue expression label. We compared the two

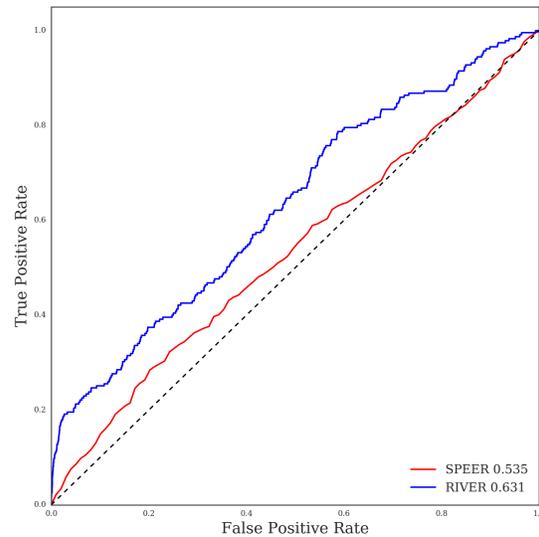


Figure 4. ROC curve comparing predictive performance of SPEER to RIVER on their respective tasks (predicting tissue-specific expression labels versus shared tissue expression labels respectively) using held-out pairs of individuals with identical variants nearby a specific gene.

274 methods at their respective tasks, predicting tissue-specific held-out expression labels
 275 and shared tissue held-out expression labels (Fig. 4). For evaluation, we followed a
 276 similar approach to Li et al. (2016) by holding out pairs of individuals that share the
 277 same rare variants nearby a specific gene. After training SPEER and RIVER on the
 278 remaining data, we computed SPEER and RIVER scores for the first individual and
 279 compared these scores to the held-out expression labels for the second individual. We
 280 observed significant performance boosts at predicting held-out shared tissue expression
 281 labels using RIVER compared to held-out tissue-specific expression labels using SPEER.
 282 These results show that tissue-specific expression labels are noisier and simply harder
 283 to predict. We investigated this further by computing the correlation between the gene
 284 expression labels across all pairs of individuals with the same rare variants. We observed
 285 a 5x increase in correlation when using shared tissue expression labels (Kendall's tau
 286 rank correlation, $\rho = 0.144$, p-value $< 1.33 \cdot 10^{-124}$) instead of tissue-specific expression
 287 labels ($\rho = 0.033$, p-value $< 3.27 \cdot 10^{-12}$).

288 **CONCLUSION.**

289 Rare variant prediction is an important problem for understanding the heritability of a
290 large number of diseases. Understanding the functional consequences of these variants
291 is a critical hurdle in our efforts towards personalized genomics. Because most diseases
292 are known to have tissue-specific molecular consequences, the development of variant
293 prediction tools that use tissue and cell-type specific context remain essential. Here we
294 have developed a probabilistic model that provides tissue-specific functional predictions
295 for rare variants. Our method shares information across tissues in order to make reliable
296 predictions.

297
298 Using our method, we observe significant performance boosts in predicting tissue-
299 specific allele-specific expression compared to the state-of-the-art, including genome
300 only prediction tools such as CADD and VEP and integrative methods like RIVER.
301 We also highlight the model's predictive power using simulated data. The simulation
302 highlights SPEER's particular usefulness with low resources across diverse tissues.

303
304 A future direction for this work is to leverage the information sharing across tissues
305 in order to make single tissue functional predictions. This will be a necessary step
306 forward given the large number of datasets with limited resources. However, predicting
307 the molecular consequences in single tissues remains a difficult problem for learning
308 reliable parameters and evaluating model performance due to noisy transcriptomic reads.

309
310 The primary application of SPEER described here involves the use of tissue-specific
311 gene expression. However, this method may also be useful for predicting alternative
312 splicing using isoform ratios or allelic imbalance using allele-specific expression by
313 direct integration of these data sources.

314

315 **Data availability** The GTEx V6 release genotypes and allele-specific expression
316 data are available on dbGaP (study accession phs000424.v6.p1; http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000424.v6.p1). GTEx V6p release expression data is available on the GTEx portal (<http://www.gtexportal.org>)

320 REFERENCES

- 321 Aguet, F., Brown, A. A., Castel, S., Davis, J. R., Mohammadi, P., Segre, A. V., Zappala,
322 Z., Abell, N. S., Fresard, L., Gamazon, E. R., Gelfand, E., Gloudemans, M. J., He, Y.,
323 Hormozdiari, F., Li, X., Li, X., Liu, B., Garrido-Martin, D., Ongen, H., Palowitch,
324 J. J., Park, Y., Peterson, C. B., Quon, G., Ripke, S., Shabalin, A. A., Shimko, T. C.,
325 Strober, B. J., Sullivan, T. J., Teran, N. A., Tsang, E. K., Zhang, H., Zhou, Y.-H.,
326 Battle, A., Bustamonte, C. D., Cox, N. J., Engelhardt, B. E., Eskin, E., Getz, G.,
327 Kellis, M., Li, G., MacArthur, D. G., Nobel, A. B., Sabbati, C., Wen, X., Wright, F. A.,
328 Lappalainen, T., Ardlie, K. G., Dermitzakis, E. T., Brown, C. D., and Montgomery,
329 S. B. (2016). Local genetic effects on gene expression across 44 human tissues.
330 *bioRxiv*.
- 331 Backenroth, D., Kiryluk, K., Xu, B., Pethukova, L., Vardarajan, B., Khurana, E., Chris-
332 tiano, A., Buxbaum, J., and Ionita-Laza, I. (2016). Tissue-specific functional effect
333 prediction of genetic variation and applications to complex trait genetics. *bioRxiv*.
- 334 Consortium, . G. P. (2015). A global reference for human genetic variation. *Nature*,
335 526(7571):68–74.
- 336 Consortium, E. P. (2012). An integrated encyclopedia of DNA elements in the human
337 genome. *Nature*, 489(7414):57–74.
- 338 Efron, B. and Tibshirani, R. J. (1994). *An introduction to the bootstrap*. CRC press.
- 339 Ernst, J. and Kellis, M. (2015). Large-scale imputation of epigenomic datasets for
340 systematic annotation of diverse human tissues. *Nat Biotech*, 33(4):364–376.
- 341 Ferrari, S. and Cribari-Neto, F. (2004). Beta regression for modelling rates and propor-
342 tions. *Journal of Applied Statistics*, 31(7):799–815.
- 343 Kircher, M., Witten, D. M., Jain, P., O’Roak, B. J., Cooper, G. M., and Shendure, J.
344 (2014). A general framework for estimating the relative pathogenicity of human
345 genetic variants. *Nature genetics*, 46(3):310–315.
- 346 Kundaje, A., Meuleman, W., Ernst, J., Bilenky, M., Yen, A., Heravi-Moussavi, A.,
347 Kheradpour, P., Zhang, Z., Wang, J., and Ziller, M. J. (2015). Integrative analysis of
348 111 reference human epigenomes. *Nature*, 518(7539):317–330.
- 349 Li, X., Kim, Y., Tsang, E. K., Davis, J. R., Damani, F. N., Chiang, C., Zappala,
350 Z., Strober, B. J., Scott, A. J., Ganna, A., Merker, J., Hall, I. M., Battle, A., and
351 Montgomery, S. B. (2016). The impact of rare variation on gene expression across
352 tissues. *bioRxiv*.
- 353 McCallum, A. K., Rosenfeld, R., Mitchell, T. M., and Ng, A. Y. (1998). Improving text
354 classification by shrinkage in a hierarchy of classes. *International Conference on*
355 *Machine Learning (ICML)*, pages 359–367.
- 356 Nelson, M. R., Wegmann, D., Ehm, M. G., Kessner, D., St. Jean, P., Verzilli, C., Shen,
357 J., Tang, Z., Bacanu, S.-A., Fraser, D., Warren, L., Aponte, J., Zawistowski, M.,
358 Liu, X., Zhang, H., Zhang, Y., Li, J., Li, Y., Li, L., Woollard, P., Topp, S., Hall,
359 M. D., Nangle, K., Wang, J., Abecasis, G., Cardon, L. R., Zöllner, S., Whittaker, J. C.,
360 Chisoe, S. L., Novembre, J., and Mooser, V. (2012). An abundance of rare functional
361 variants in 202 drug target genes sequenced in 14,002 people. *Science (New York,*
362 *N.Y.)*, 337(6090):100–104.
- 363 Quang, D., Chen, Y., and Xie, X. (2015). DANN: A deep learning approach for
364 annotating the pathogenicity of genetic variants. *Bioinformatics*, 31(5):761–763.

- 365 Raina, R., Ng, A., and Koller, D. (2006). Constructing informative priors using transfer
366 learning. *Proceedings of the 23rd international conference on Machine learning*, page
367 713–720.
- 368 Stegle, O., Parts, L., Piipari, M., Winn, J., and Durbin, R. (2012). Using probabilistic es-
369 timation of expression residuals (PEER) to obtain increased power and interpretability
370 of gene expression analyses. *Nature protocols*, 7(3):500–7.
- 371 Tennessen, J. A., Bigham, A. W., O’Connor, T. D., Fu, W., Kenny, E. E., Gravel, S.,
372 McGee, S., Do, R., Liu, X., Jun, G., Kang, H. M., Jordan, D., Leal, S. M., Gabriel, S.,
373 Rieder, M. J., Abecasis, G., Altshuler, D., Nickerson, D. A., Boerwinkle, E., Sunyaev,
374 S., Bustamante, C. D., Bamshad, M. J., Akey, J. M., GO, B., GO, S., and Project, o. b.
375 o. t. N. E. S. (2012). Evolution and Functional Impact of Rare Coding Variation from
376 Deep Sequencing of Human Exomes. *Science*, 337(6090):64–69.
- 377 Thrun, S. (1996). Is learning the n-th thing any easier than learning the first? *Advances*
378 *in neural information processing systems*, pages 640–646.
- 379 Uricchio, L. H., Zaitlen, N. A., Ye, C. J., Witte, J. S., and Hernandez, R. D. (2016).
380 Selection and explosive growth alter genetic architecture and hamper the detection of
381 causal rare variants. *Genome Research*, 26(7):863–873.
- 382 Yan, H., Yuan, W., Velculescu, V. E., Vogelstein, B., and Kinzler, K. W. (2002). Allelic
383 variation in human gene expression. *Science*, 297(5584):1143.
- 384 Zhang, K., Li, J. B., Gao, Y., Egli, D., Xie, B., Deng, J., Li, Z., Lee, J.-H., Aach, J.,
385 and Leproust, E. M. (2009). Digital RNA allelotyping reveals tissue-specific and
386 allele-specific gene expression in human. *Nature methods*, 6(8):613–618.

387 **SUPPLEMENT.**

388 **Tissue groups.** We evaluated the pairwise similarity between gene expression patterns
389 across tissues and identified the following list of tissue groups used in the GTEx results
390 section. GTEx ids are listed below:

391 **Brain** Brain Caudate basal ganglia, Brain Nucleus accumbens basal ganglia, Brain
392 Putamen basal ganglia, Brain Anterior cingulate cortex BA24, Brain Cortex, Brain
393 Frontal Cortex BA9

394 **Artery and Fat** Artery Coronary, Artery Aorta, Artery Tibial, Esophagus Muscularis,
395 Esophagus Gastroesophageal Junction, Colon Sigmoid, Adipose Subcutaneous,
396 Adipose Visceral Omentum, Breast Mammary Tissue

397 **Muscle** Muscle Skeletal, Heart Atrial Appendage, Heart Left Ventricle

398 **Epithelial** Skin Not Sun Exposed Suprapubic, Skin Sun Exposed Lower leg, Esophagus
399 Mucosa, Vagina

400 **Digestive** Stomach, Colon Transverse, Lung, Thyroid, Prostate

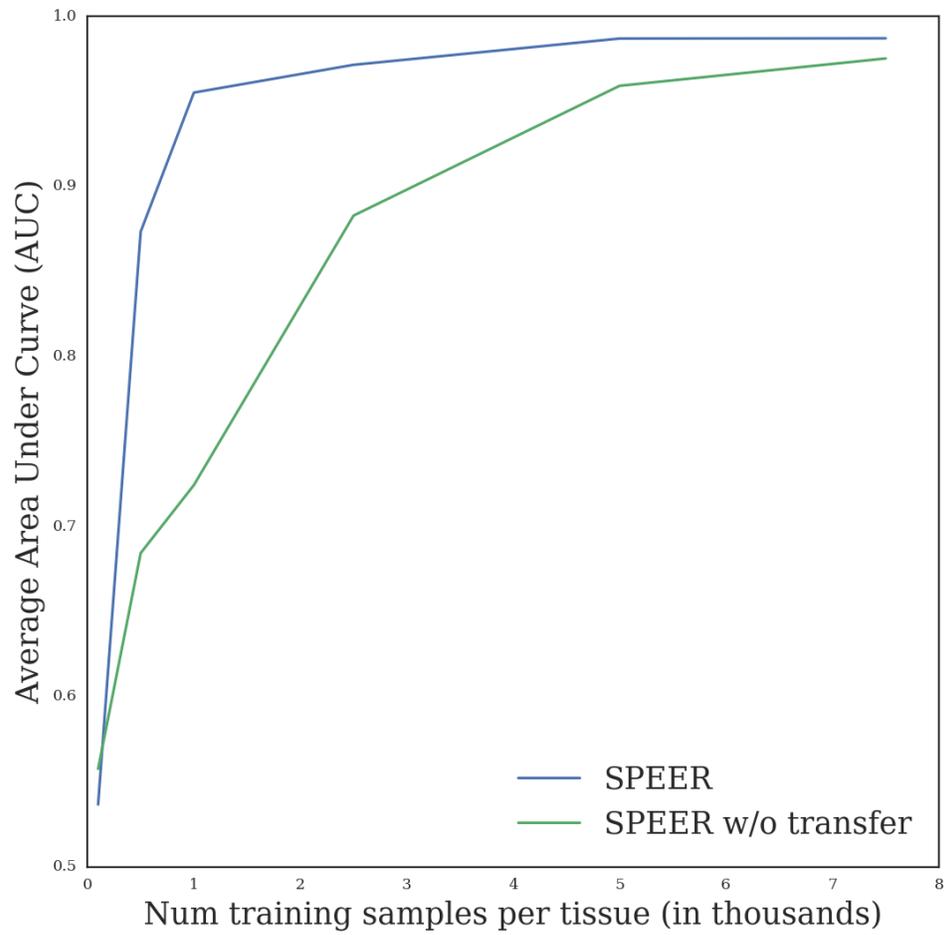


Figure 5. Area under curve (AUC) averaged across five tissue groups for different number of training samples using simulated data.

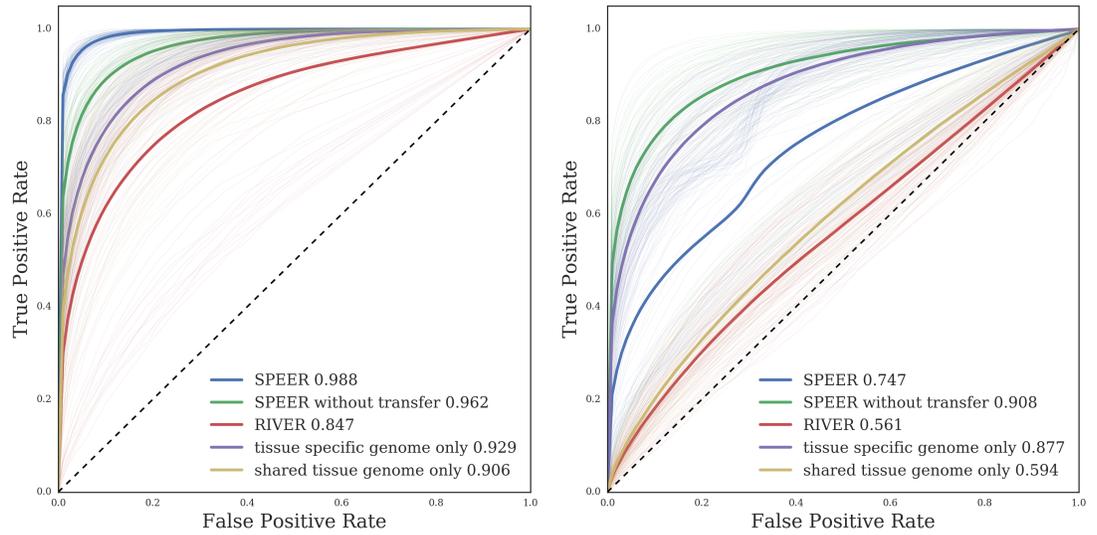


Figure 6. ROC curves for equal effects setting comparing SPEER to four benchmarks in the tied tissue simulation (left) and the independent tissue simulation (right).

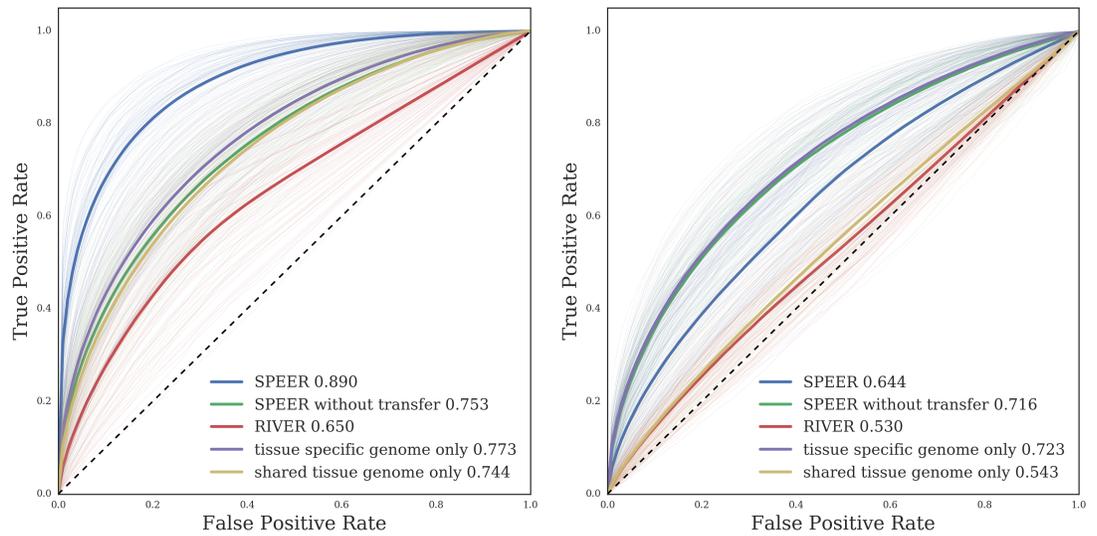


Figure 7. ROC curves for weaker effects setting comparing SPEER to four benchmarks in the tied tissue simulation (left) and the independent tissue simulation (right).

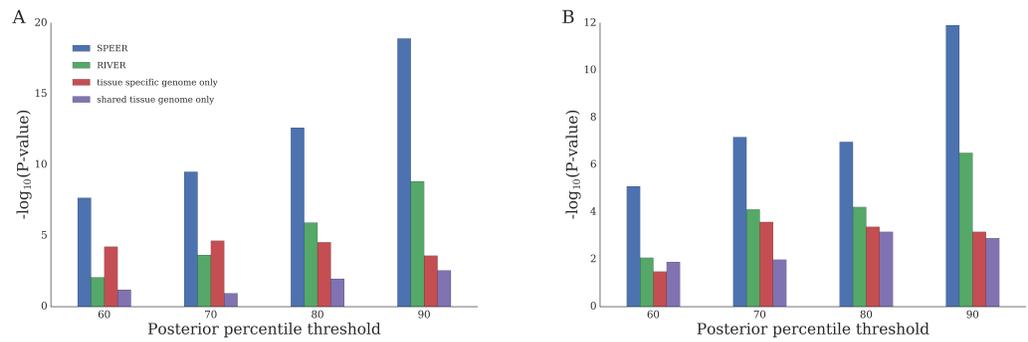


Figure 8. SPEER scores compared to tissue-specific allelic imbalance using z-score expression outlier thresholds of 1.75 (left) and 2.0 (right) in GTEx data.

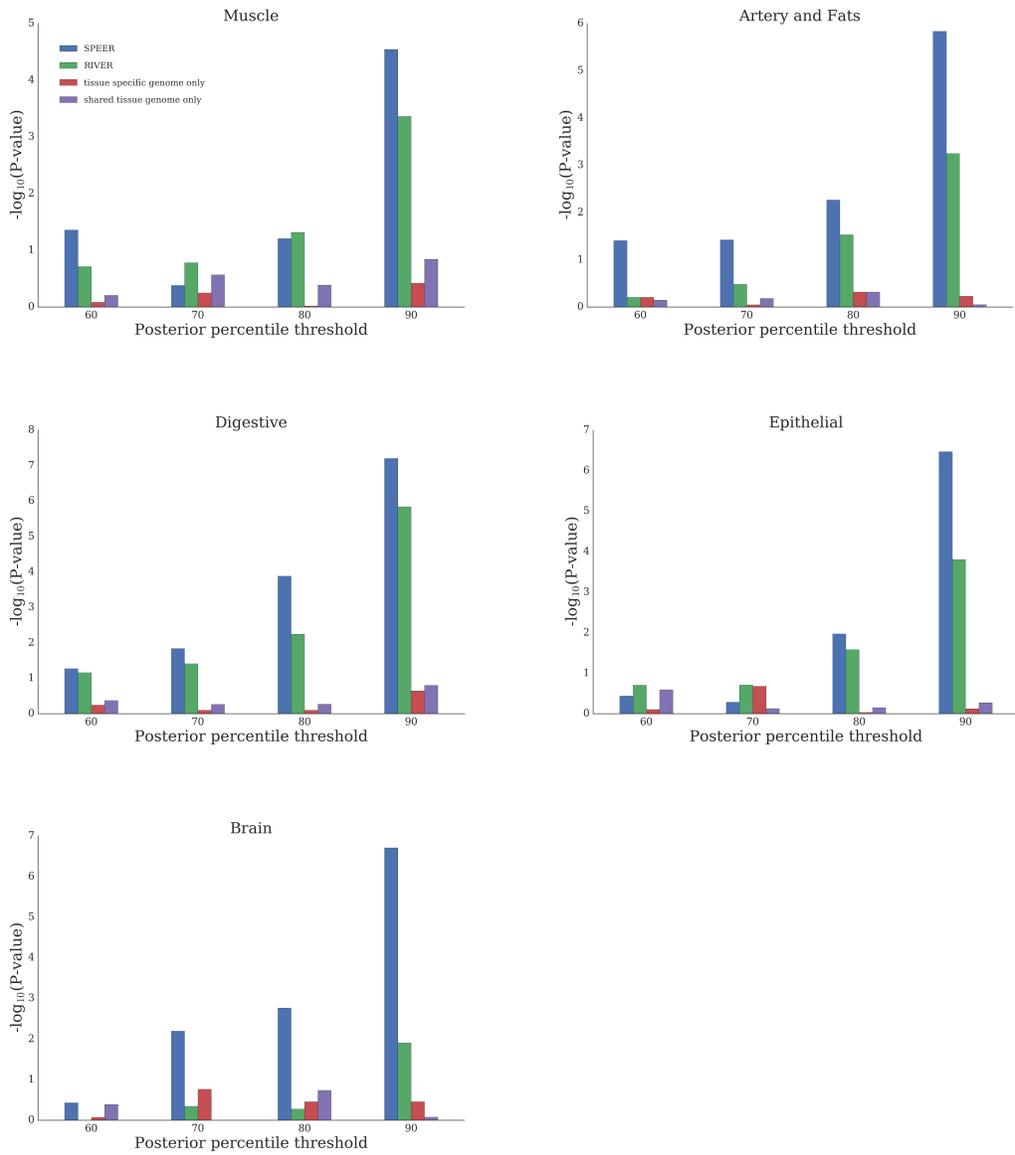


Figure 9. SPEER scores compared to allelic imbalance in the five tissue groups using GTEx data.